

Abstract Interpretation

(Cousot, Cousot 1977)

also known as

Data-Flow Analysis

(Kildall 1973)

int a, b, step, i;

boolean c;

a = 0;

b = a + 10;

step = -1;

if (step > 0) {

i = a;

} else {

i = b;

}

c = true;

while (c) {

print(i);

i = i + step; // can emit decrement

if (step > 0) {

c = (i < b);

} else {

c = (i > a); // can emit better instruction here

} // insert here (a = a + step), redo analysis

}

Example: Constant propagation

Here is why it is useful

$a \rightarrow \perp$ $b \rightarrow \perp$
 $step \rightarrow \perp$
 $i \rightarrow \perp = \phi$

$a \rightarrow \{0\}$
 $b \rightarrow \{10\}$
 $step \rightarrow \{-1\}$

$i \rightarrow T$

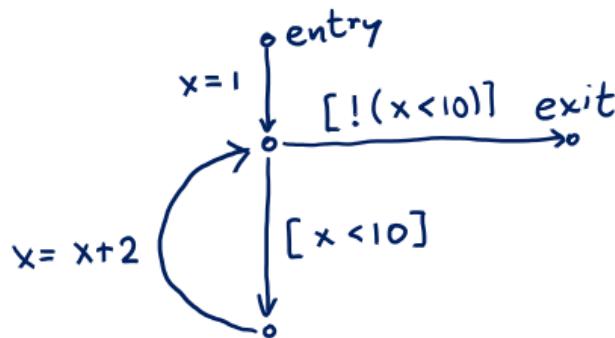
Goal of Data-Flow Analysis

Automatically compute information about the program

- Use it to report errors to user (like type errors)
- Use it to optimize the program

Works on control-flow graphs:
(like flow-charts)

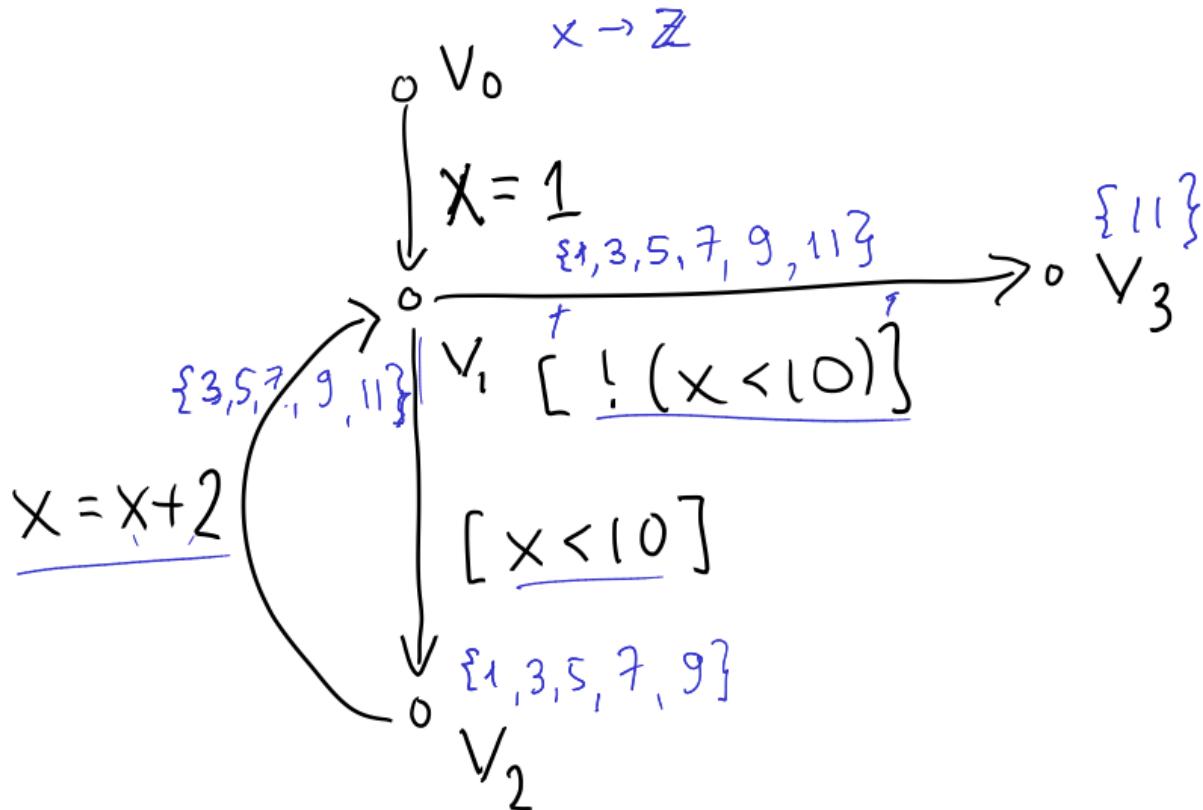
```
x = 1
while (x < 10) {
    x = x + 2
}
```



Interpretation and Abstract Interpretation

- Control-Flow graph is similar to AST
- We can
 - interpret control flow graph
 - generate machine code from it (e.g. LLVM, gcc)
 - abstractly interpret it: do not push values, but
approximately compute supersets of possible values
(e.g. intervals, types, etc.)

Compute Range of x at Each Point $\in \mathbb{Z}$



What we see in the sequel

1. How to compile abstract syntax trees into control-flow graphs
2. Lattices, as structures that describe abstractly sets of program states (facts)
3. Transfer functions that describe how to update facts
4. Basic idea of fixed-point iteration

Generating Control-Flow Graphs

- Start with graph that has one entry and one exit node and label is entire program
- Recursively decompose the program to have more edges with simpler labels
- When labels cannot be decomposed further, we are done

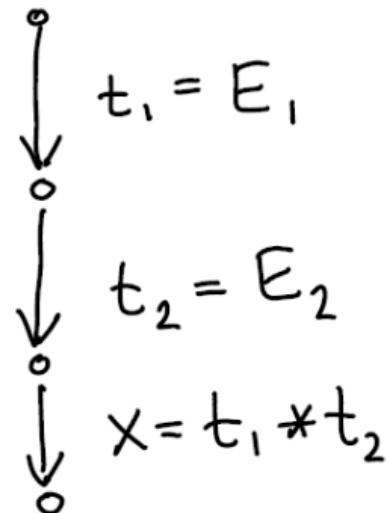
Flattening Expressions for simplicity and ordering of side effects

E_1, E_2 - complex expressions

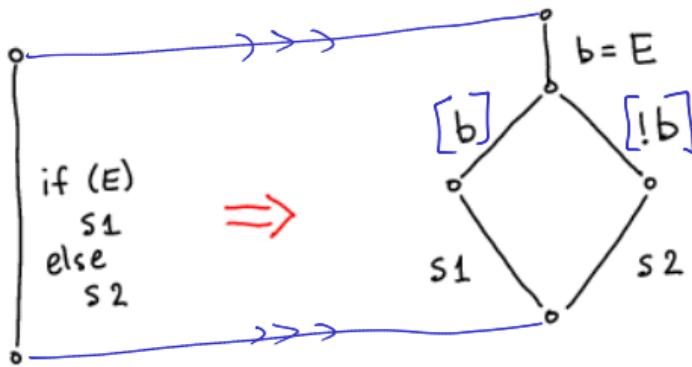
t_1, t_2 - fresh variables



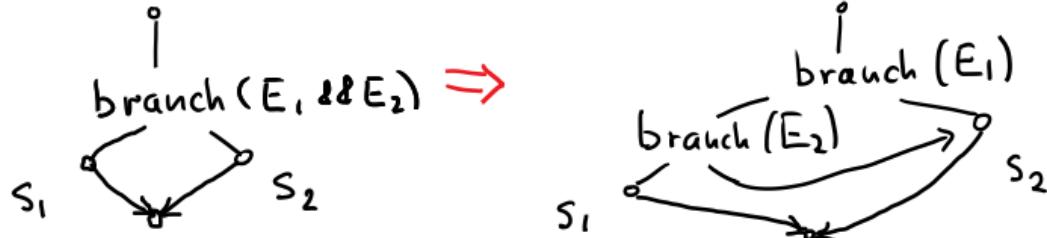
$$x = E_1 * E_2 \quad \Rightarrow$$



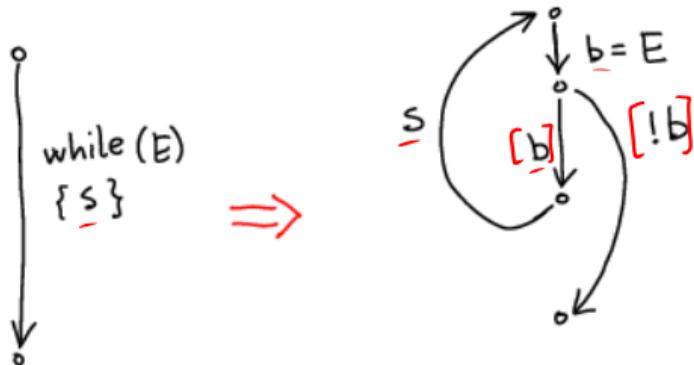
If-Then-Else



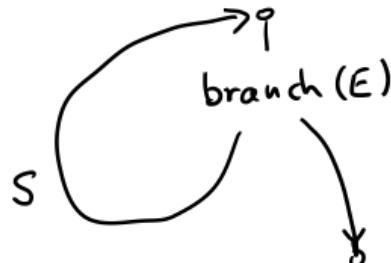
Translation using "branch" instruction works nicely for control flow graphs: two destinations



While



Better translation uses the "branch" instruction



Example 1: Convert to CFG

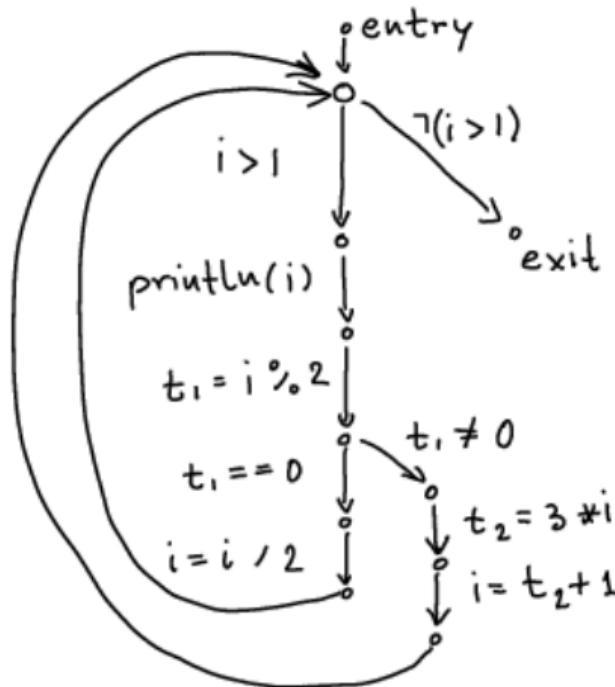
```
while (i < 10) {  
    println(j);  
    i = i + 1;  
    j = j +2*i + 1;  
}
```

Example 2: Convert to CFG

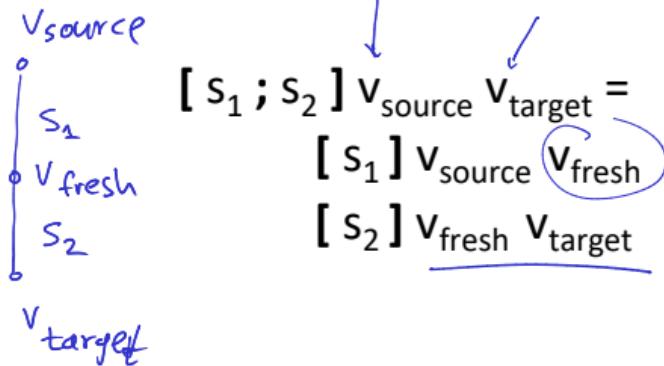
```
int i = n;  
while (i > 1) {  
    println(i);  
    if (i % 2 == 0) {  
        i = i / 2;  
    } else {  
        i = 3*i + 1;  
    }  
}
```

Example 2 Result

```
int i = n;  
while (i > 1) {  
    println(i);  
    if (i % 2 == 0) {  
        i = i / 2;  
    } else {  
        i = 3*i + 1;  
    }  
}
```



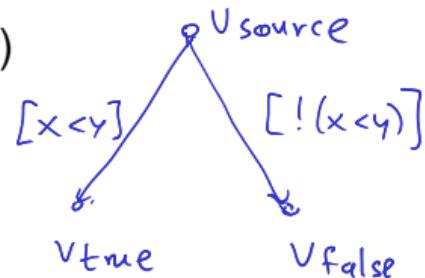
Translation Functions



insert (v_s ,stmt, v_t)=
 $cfg = cfg + (v_s, \text{stmt}, v_t)$

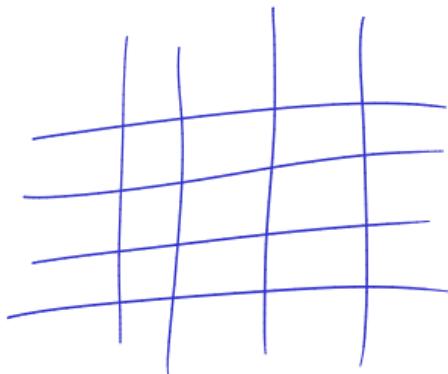
$[x=y+z] v_s v_t = \text{insert}(v_s, x=y+z, v_t)$

branch($x < y$)

$$v_{source} v_{true} v_{false} = \text{insert}(v_{source}, [x < y], v_{true});$$
$$\text{insert}(v_{source}, [!(x < y)], v_{false})$$


when y,z are constants or variables

Abstract Interpretation: Analysis Domain (D) Lattices



Lattice

\subseteq

Partial order: binary relation \leq (subset of some D^2)

which is

- reflexive: $x \leq x$
- anti-symmetric: $x \leq y \wedge y \leq x \rightarrow x = y$
- transitive: $x \leq y \wedge y \leq z \rightarrow x \leq z$

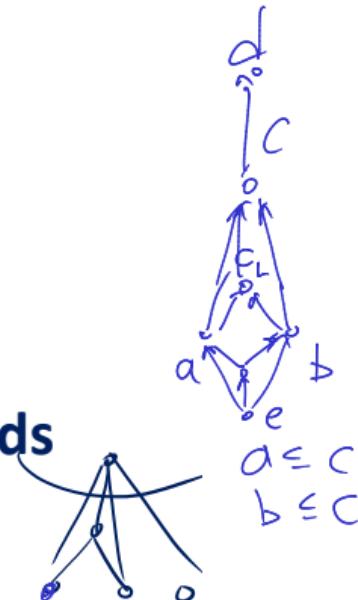
closure of
directed
acyclic
acyclic

Lattice is a partial order in which every two-element set has **least among its upper bounds** and **greatest among its lower bounds**

- Lemma: if (D, \leq) is lattice and D is finite, then lub and glb exist for every finite set

\sqcap \sqcup

$\sqcup \{a, b, c\}$

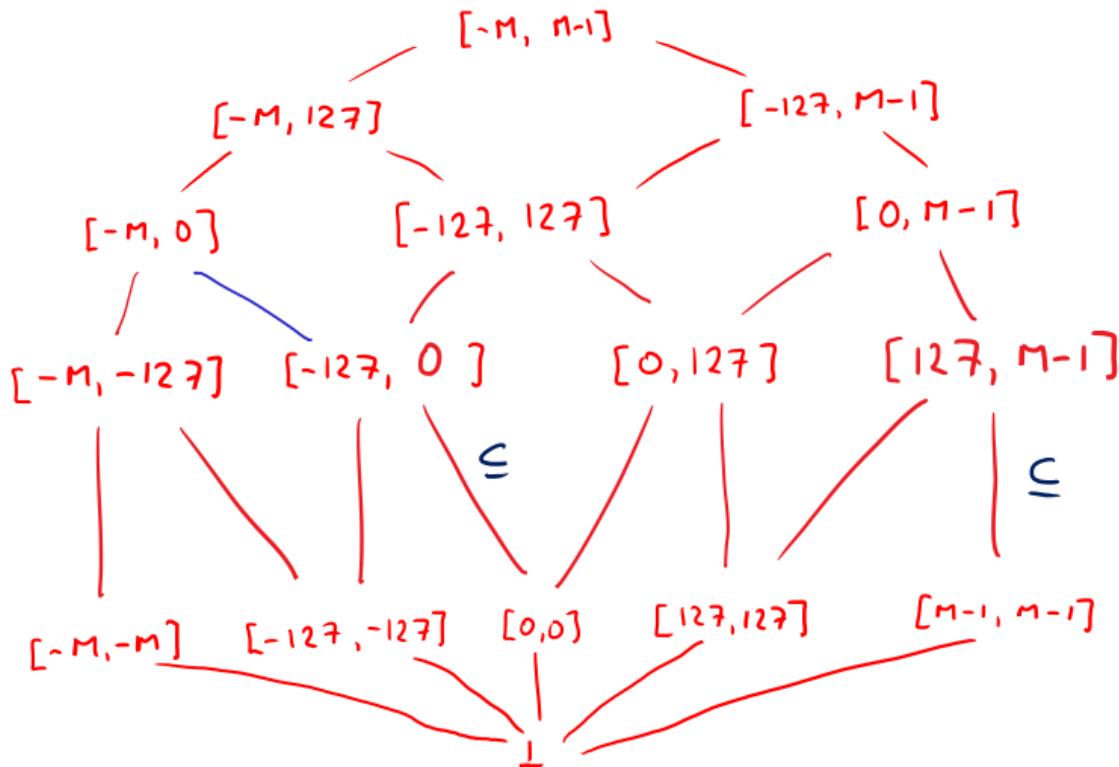


Graphs and Partial Orders

- If the domain is finite, then partial order can be represented by directed graphs
 - if $x \leq y$ then draw edge from x to y
- For partial order, no need to draw $x \leq z$ if $x \leq y$ and $y \leq z$. So we only draw non-transitive edges
- Also, because always $x \leq x$, we do not draw those self loops
- Note that the resulting graph is acyclic: if we had a cycle, the elements must be equal

Domain of Intervals $[a,b]$ where $a,b \in \{-M, -127, 0, 127, M-1\}$

$$M = 2^{31}$$



Defining Abstract Interpretation

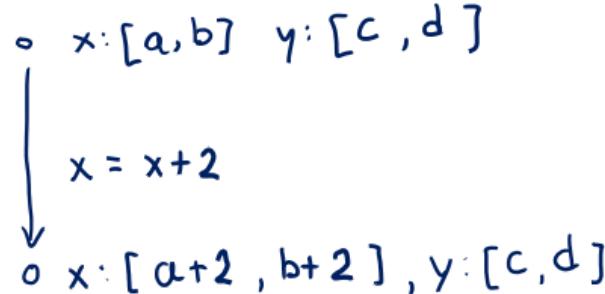
Abstract Domain D describing which information to compute – this is often a lattice

- inferred types for each variable: $x:T_1, y:T_2$
- interval for each variable $x:[a,b], y:[a',b']$

Transfer Functions, $[[st]]$ for each statement st , how this statement affects the facts $D \rightarrow D$

- Example:

$$\begin{aligned} [[x = x+2]](x:[a,b], \dots) \\ = (x:[a+2, b+2], \dots) \end{aligned}$$



For now, we consider arbitrary integer bounds for intervals

- Thus, we work with BigInt-s
- Often we must analyze machine integers
 - need to correctly represent (and/or warn about) overflows and underflows
 - fundamentally same approach as for unbounded integers
- For efficiency, many analysis do not consider arbitrary intervals, but only a subset of them W
- We consider as the domain
 - empty set (denoted \perp , pronounced “bottom”)
 - all intervals $[a,b]$ where a,b are integers and $a \leq b$, or where we allow $a = -\infty$ and/or $b = \infty$
 - set of all integers $[-\infty, \infty]$ is denoted T, pronounced “top”

$[1, 9]$
 $1, 3, 5, 7, 9$

Find Transfer Function: Plus

Suppose we have only two integer variables: x, y

$$\begin{array}{l} \bullet \quad x: [a,b] \quad y: [c,d] \\ \downarrow \\ \bullet \quad x = x+y \\ \quad x: [a',b'] \quad y: [c',d'] \end{array}$$

If $a \leq x \leq b$ $c \leq y \leq d$
and we execute $x = x+y$
then $x' = x+y$ $x = x'-y$
 $y' = y$ $a \leq x \leq b$
 $a \leq x' - y \leq b$
so $a+c \leq x' \leq b+d$ $a+c \leq x+y \leq b+d$
 $c \leq y' \leq d$

So we can let

$$\begin{array}{ll} a' = a+c & b' = b+d \\ c' = c & d' = d \end{array}$$

Find Transfer Function: Minus

Suppose we have only two integer variables: x, y

$$\begin{array}{l} x : [a, b] \quad y : [c, d] \\ \downarrow \\ y = x - y \\ \downarrow \\ x' : [a', b'] \quad y' : [c', d'] \end{array}$$

If

and we execute $y = x - y$

then

$$y' = x - y$$

$$y' = x - y'$$

$$c \leq y \leq d$$

$$c \leq x - y' \leq d$$

$$c \leq x - y'$$

$$c \leq x - c \leq b - c$$

So we can let

$$a' = a \quad b' = b$$

$$c' = a - d \quad d' = b - c$$

Transfer Functions for Tests

Tests e.g. $[x > 1]$ come from translating if,while into
CFG

$$x : [-10, 10]$$

if ($x > 1$) {

$$x : [2, 10]$$

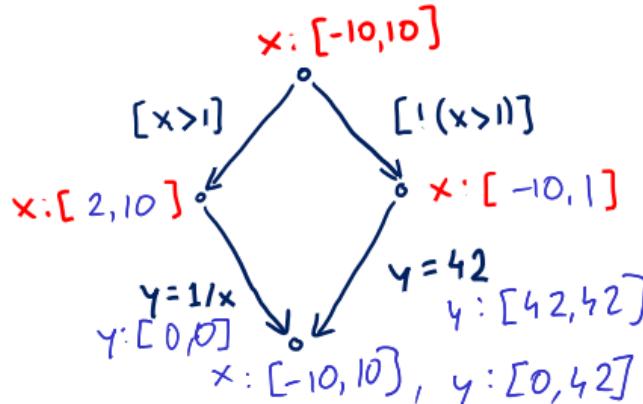
$$y = 1/x$$

} else {

$$x : [-10, 1]$$

$$y = 42$$

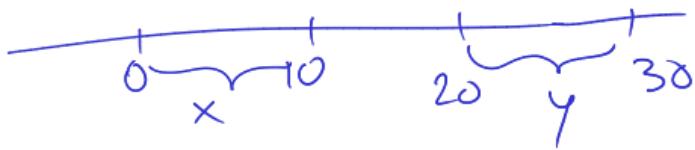
}



$$x : [a, b] \quad y : [c, d]$$

$$[x > y]$$

\perp



Joining Data-Flow Facts

$x: [-10, 10] \quad y: [-1000, 1000]$

if ($x > 0$) {

$x: [1, 10] \quad y: [-1000, 1000]$

$y = x + 100$

$x: [1, 10] \quad y: [101, 110]$

} else { $x \leq 0$

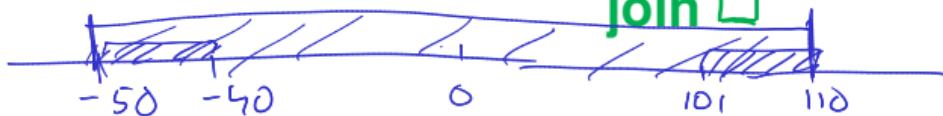
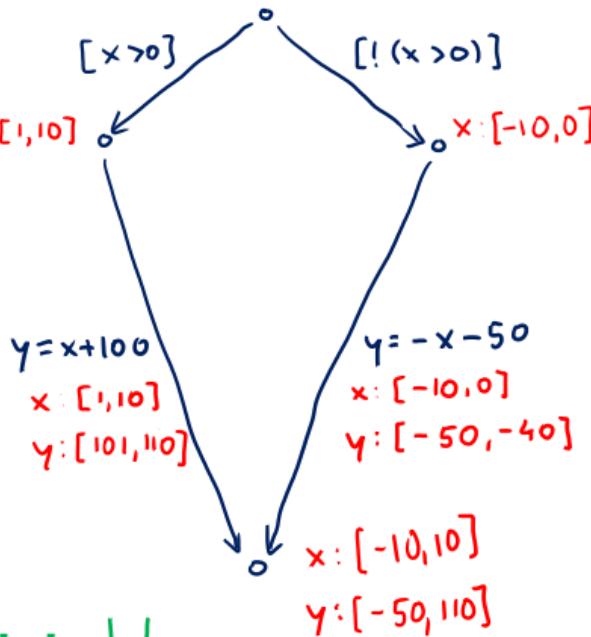
$x: [-10, 0] \quad y: [-1000, 1000]$

$y = -x - 50$

$x: [-10, 10] \quad y: [-50, -40]$

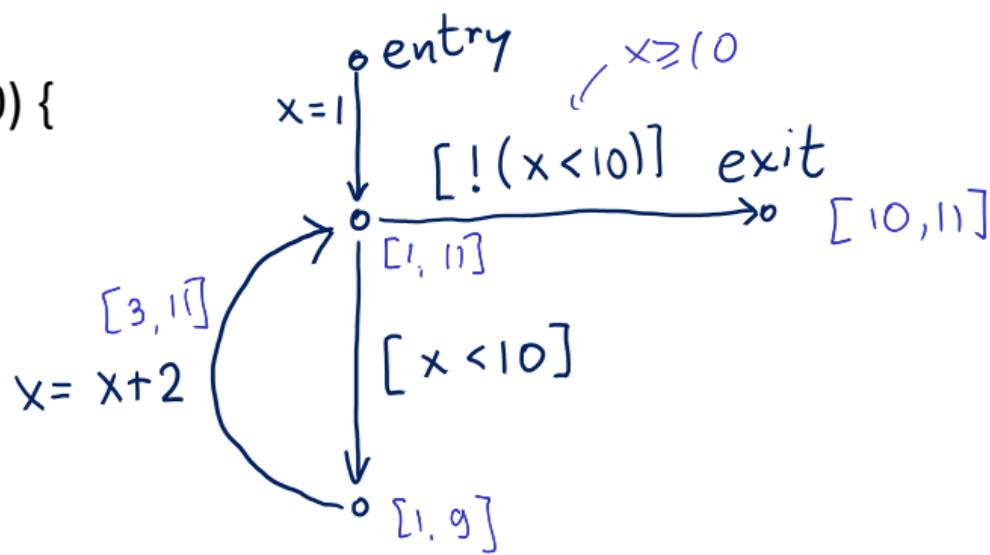
}

$x: \quad y:$



Handling Loops: Iterate Until Stabilizes

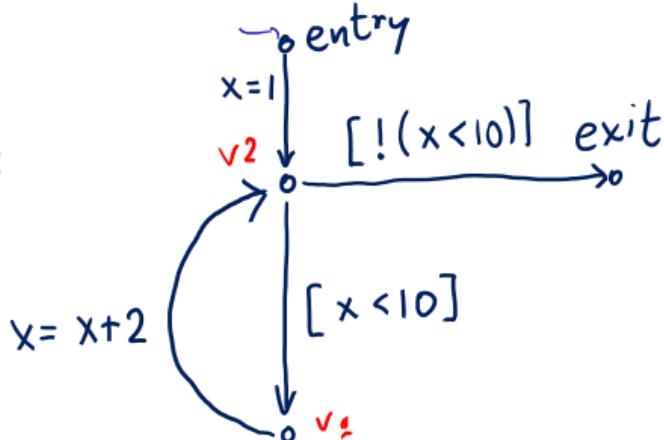
```
x = 1  
[1,1]  
while(x < 10) {  
    [1,1]  
    x = x + 2  
    [3,3]  
}
```



Analysis Algorithm

```
var facts : Map[Node,Domain] = Map.withDefault(empty)
facts(entry) = initialValues
while (there was change)
    pick edge (v1,stmt,v2) from CFG
        such that facts(v1) has changed
        facts(v2)=facts(v2) join transferFun(stmt, facts(v1))
}
```

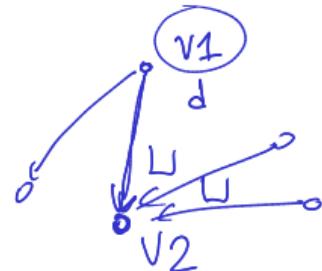
Order does not matter for the end result, as long as we do not permanently neglect any edge whose source was changed.



```

→ var facts : Map[Node,Domain] = Map.withDefault(empty)
    var worklist : Queue[Node] = empty
    {def assign(v1:Node,d:Domain) = if (facts(v1)!=d) {
        facts(v1)=d
        for (stmt,v2) <- outEdges(v1) { worklist.add(v2) }
    }
    assign(entry, initialValues)
    while (!worklist.isEmpty) {
        var v2 = worklist.getAndRemoveFirst
        var update = facts(v2)
        for (v1,stmt) <- inEdges(v2)
            { update = update join transferFun(facts(v1),stmt) }
        assign(v2, update)
    }
}

```



Work List Version

Exercise: Run range analysis, prove that **error** is unreachable

```
int M = 16;  
int[M] a;  
x := 0;  
while (x < 10) {  
    x := x + 3;  
}  
if (x >= 0) {  
    if (x <= 15)  
        a[x]=7;  
    else  
        error;  
} else {  
    error;  
}
```

checks array accesses

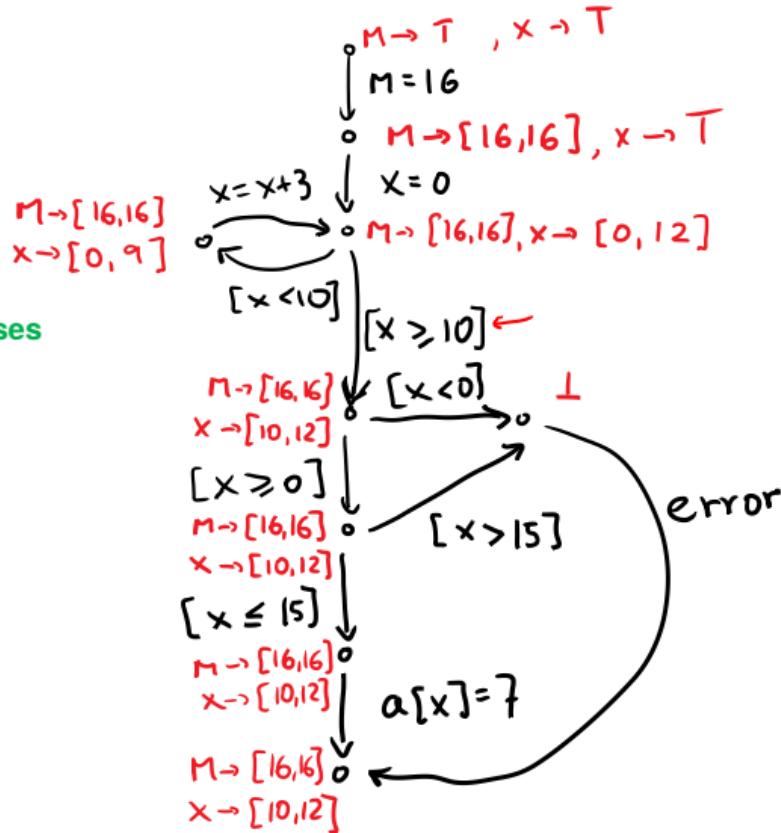


Range analysis results

```

int M = 16;
int[M] a;
x := 0;
while (x < 10) {
    x := x + 3;
}
if (x >= 0) {
    if (x <= 15)
        a[x]=7;
    else
        error;
} else {
    error;
}
    
```

checks array accesses

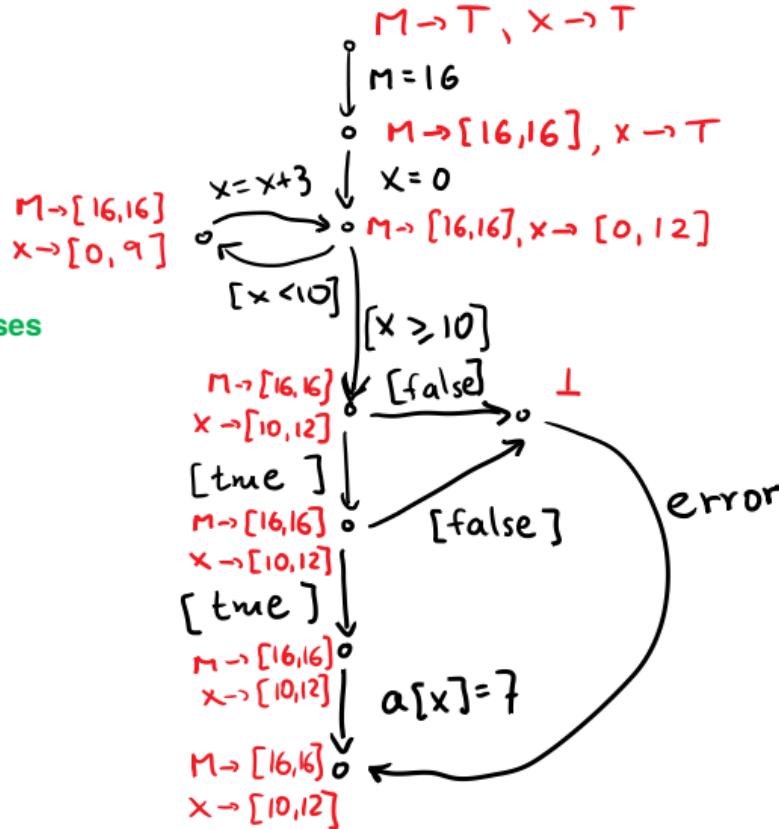


Simplified Conditions

```

int M = 16;
int[M] a;
x := 0;
while (x < 10) {
    x := x + 3;
}
if (x >= 0) {
    if (x <= 15)
        a[x]=7;
    else
        error;
} else {
    error;
}
    
```

checks array accesses



Remove Trivial Edges, Unreachable Nodes

```
int M = 16;
```

```
int[M] a;
```

```
x := 0;
```

```
while (x < 10) {
```

```
    x := x + 3;
```

```
}
```

checks array accesses

```
if (x >= 0) {
```

```
    if (x <= 15)
```

```
        a[x]=7;
```

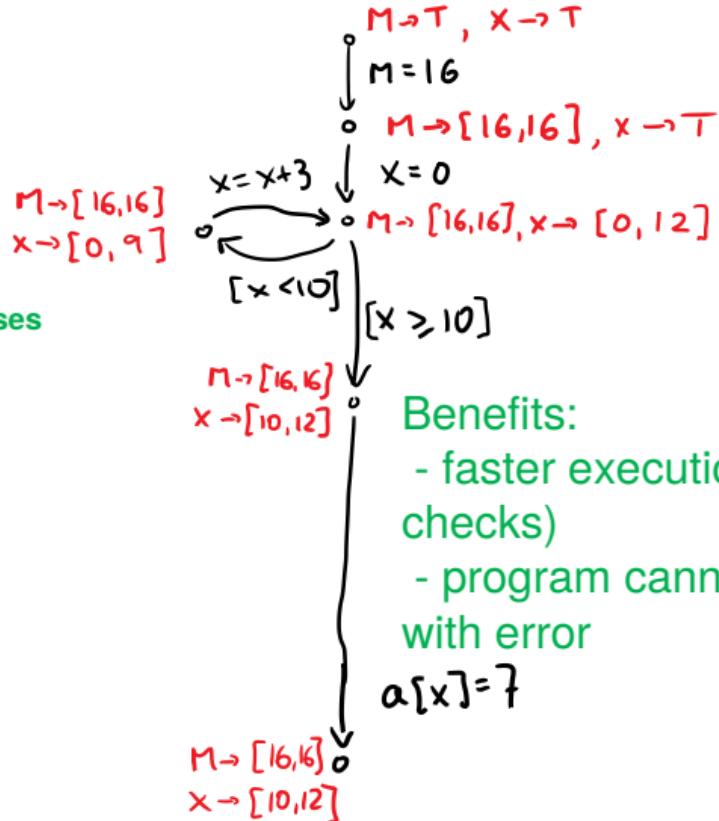
```
    else
```

```
        error;
```

```
} else {
```

```
    error;
```

```
}
```



Benefits:

- faster execution (no checks)
- program cannot crash with error

Constant Propagation Domain

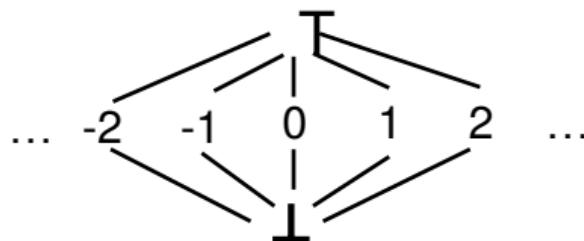
Domain values D are:

- intervals $[a,a]$, denoted simply ‘ a ’
- empty set, denoted \perp and set of all integers T

Formally, if \mathbb{Z} denotes integers, then

$$D = \{\perp, T\} \cup \{ a \mid a \in \mathbb{Z}\}$$

D is an infinite set



Constant Propagation Transfer Functions

$$\begin{array}{c} v_1 \\ \downarrow \\ x = y + z \\ \downarrow \\ v_2 \end{array}$$

$$\begin{aligned} S_x + S_y \\ = \{x+y \mid x \in S_x, y \in S_y\} \end{aligned}$$

For each variable (x,y,z) and each CFG node (program point) we store: \perp , a constant, or T

table for $+$:

\perp	\perp	C_y	T
\perp	\perp	\perp	\perp
C_z	\perp	$C_y + C_z$	T
T	\perp	T	T

abstract class Element
case class Top extends Element
case class Bot extends Element
case class Const(v:Int) extends Element
var facts : Map[Nodes,Map[VarNames,Element]]
what executes during analysis of $x=y+z$:
oldY = facts(v₁)("y")
oldZ = facts(v₁)("z")
newX = tableForPlus(oldY, oldZ)
facts(v₂) = facts(v₂) join facts(v₁).updated("x", newX)}

```
def tableForPlus(y:Element, z:Element)
= (x,y) match {
  case (Const(cy),Const(cz)) =>
    Const(cy+cz)
  case (Bot,_) => Bot
  case (_,Bot) => Bot
  case (Top,Const(cz)) => Top
  case (Const(cy),Top) => Top
}
```